



Protecting Shared Data

Today...

- Project Assignments
 - User Study and Consent
 - Exercise (from Lecture 3)
- Protecting Publicly Shared Data
 - Social Media (e.g. Facebook, Twitter, ...)
- Publicly Shared Data (e.g. Census Data, Research Data,...)
- Network Anonymisation (e.g. Tor...)

Project Assignments...

- Some points to remember ...
 - Read to help guide project focus
 - Good start point: Project description references
 - Remember to Document! (Portfolio)

Project Assignments...

- Some approaches to conducting your project...
 1. Study an aspect of an existing system
 - Why do users post private pictures on Social media?
 2. Build a system and test
 - How to users respond to X and Y added features?
 3. Or Prototype (e.g. Paper Diagrams / Simulations /Artefacts), get data, then build
 - Will people be interested in using this?

Project Assignments...

- Other Considerations...
 - Remember to also consider adversarial scenarios
 - What can users do wrong?
 - How can data be exposed? Why? ...
 - Important:
 - Potential negative consequences of the solution
 - *"Should I build this system, just because I can?"*
- User consent is important...
 - Make sure to obtain before you collect **ANY** data whether this is in-person or electronically

Exercise #1

- Q1:
 - *Survey*
 - *Diary Study*
 - *Interviews*
 - *A Usability Test – e.g. based on a prototype*
 - *Collecting observational or experimental data in the field*
- Q2:
 - Decide which method you want to use
 - Do Q2 a., b., and c

Exercise #2 (User Study Protocol)

- Review Proposed User Study Protocol...
- Some ethics concerns:
 - Intrusive
 - Health concerns
 - Privacy of users (watching the users without informed consent)...
- Some Suggested Modifications to the Protocol:
 - User consent form
 - Perhaps more compensation for users
 - Preliminary testing of device (lab) to identify health risks...

Protecting Publicly Shared Data

SOCIAL MEDIA



“I read my Twitter the next morning and was astonished” A Conversational Perspective on twitter Regrets

Sleeper et al. (CHI 2013)

I have no idea how to send a private message...

Oh dear! That was meant to be a direct message, not a group message!

The Twitter Typo That Exposed Anthony Weiner...

Privacy on Social Media – User View

- “Context Collapse” combines separate offline groups (e.g. friends, family, coworkers...)
- "Imagined audience" may not align with actual audience
- "Temporal Changes" in preferences
- “Privacy Tools” may be unclear / hard to use (prone to failure when network conditions are not optimal)

Privacy on Social Media – User View

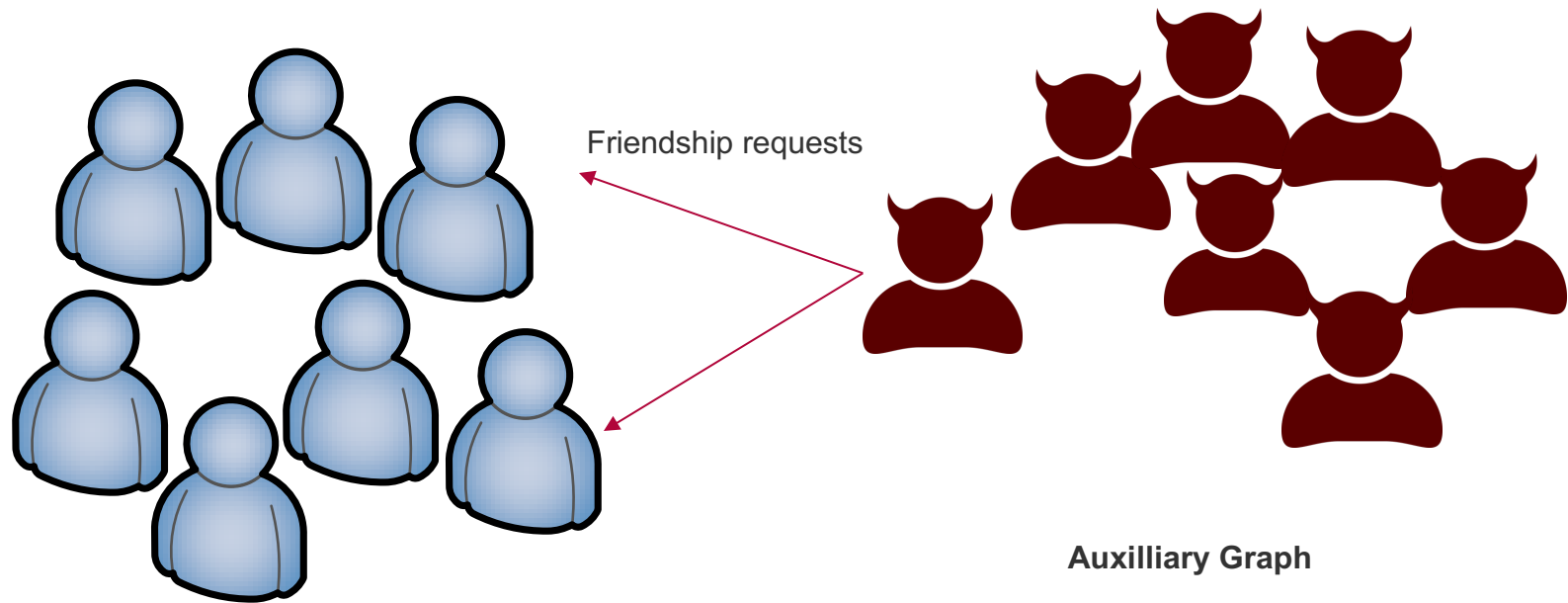
- Minor to severe consequences – due to revealed private information.
 - e.g. Regret, Loss of Credibility, difficult to delete...
- Variety of social media platforms, functionalities, focus,...
- Different types of users...
- Range of privacy threats

Privacy on Social Media – Expert View

- Social media platforms are not designed for secret sharing...
- Users should pay more attention to what they post...
- Users should learn to use privacy settings better...
- At no other time has it become more important than now to "Look Before You Click..."

De-anonymising Social Graphs

De-Anonymizing Social Graphs via Node Similarity, Fu et al. (2014)
De-Anonymizing Web Browsing Data with Social Networks, Su et al. (2017)



Target Graph
(anonymised and published social graph)

Auxilliary Graph

Attack based on prior knowledge to disclose identities

De-anonymising Social Graphs

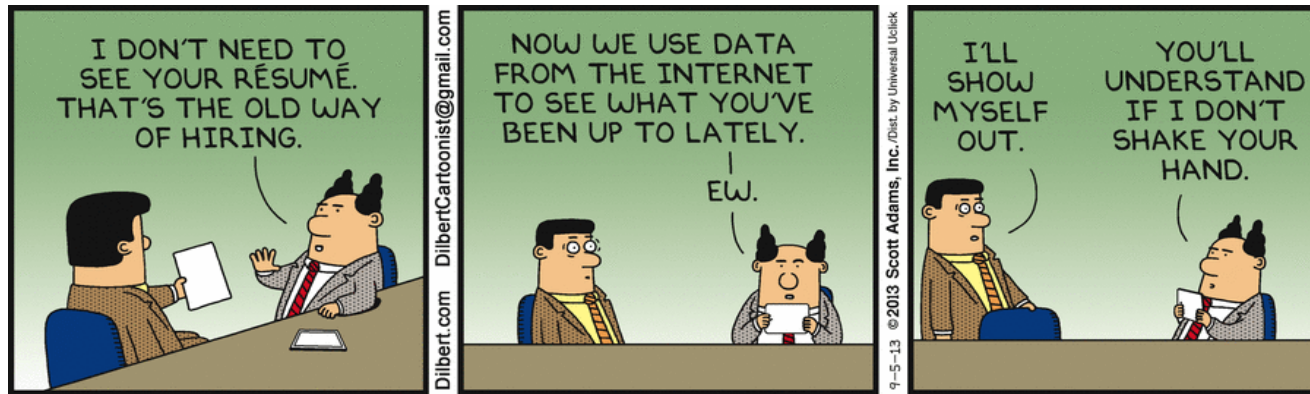
- Browsing histories can be linked to social media profiles
- **Observation:** Users are more likely to click on links posted by users whom they follow (are “friends to”)
- These patterns persist in their browsing history
- **Result:** By finding the social media profile that matches a browsing history we can de-anonymise the profile owner (user)

Public Data

- Health-care datasets
 - Clinical studies, hospital discharge databases ...
- Genetic datasets
 - \$1000 genome, HapMap, deCode ...
- Demographic datasets
 - Census Bureau, sociology studies ...
- Search logs, recommender systems, social networks, blogs ...
 - AOL search data, social networks of blogging sites, Netflix movie ratings, Amazon ...

Privacy Concerns

- How do we protect privacy here?



- Data comes from usage patterns
 - Privacy exposure is not necessarily the user's fault
 - So, how can we protect user privacy?

Anonymize the Data...

- **Syntactic** versus Semantic Anonymization

Personal Identifiers		Quasi-identifier Attributes			Sensitive Attributes
First Name	Surname	Age	Post Code	Gender	Symptoms
Alice	McCarthy	19	02138	Female	Ovarian Cancer
Steve	Joe	20	02139	Male	Flu
Eve	Boxwell	25	02141	Female	Ovarian Cancer
John	Mark	22	02142	Male	Lung cancer
Rochelle	Alers	35	02138	Female	Thyroid Disease
Ben	Colbert	30	02139	Male	Thyroid Disease

- Is this enough?

Harvard Professor Re-Identifies Volunteers in a DNA study...

Personal Identifiers		Quasi-identifier Attributes			Sensitive Attributes
First Name	Surname	Age	Post Code	Gender	Symptoms
		19	02138	Female	Ovarian Cancer
		20	02139	Male	Flu
		25	02141	Female	Ovarian Cancer
		22	02142	Male	Lung cancer
		35	02138	Female	Thyroid Disease
		30	02139	Male	Thyroid Disease

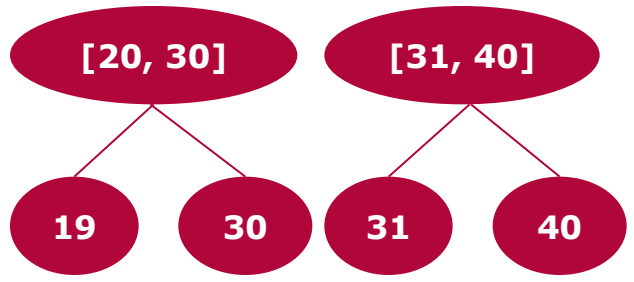
Anonymised Data

First Name	Surname	Age	Post Code	Gender
Alice	McCarthy	19	02138	Female
Steve	Joe	20	02139	Male
Eve	Boxwell	25	02141	Female
John	Mark	22	02142	Male
Rochelle	Alers	35	02138	Female
Ben	Colbert	30	02139	Male

Voter Registration Data

Generalising Attributes and k-Anonymity...

Generalisation Tree (Example)



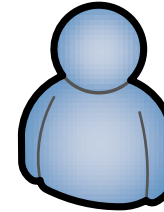
Age	Post Code	Gender	Symptoms
19	02138	Female	Ovarian Cancer
20	02139	Male	Flu
25	02141	Female	Ovarian Cancer
22	02142	Male	Lung cancer
35	02138	Female	Thyroid Disease
30	02139	Male	Thyroid Disease

Age	Post Code	Gender	Symptoms
[20, 30]	02130	Female	Ovarian Cancer
[20, 30]	02130	Male	Flu
[20, 30]	02140	Female	Ovarian Cancer
[20, 30]	02140	Male	Lung cancer
[31, 40]	02130	Female	Thyroid Disease
[31, 40]	02130	Male	Thyroid Disease

Classify in buckets of size k

Re-Identification By Linking...

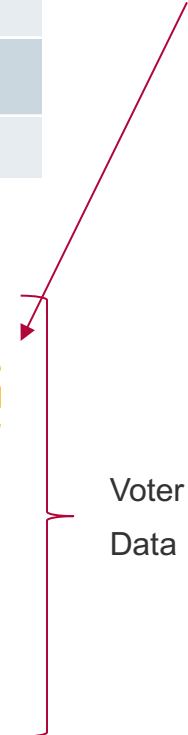
Age	Post Code	Gender	Symptoms
[20, 30]	02130	Female	Ovarian Cancer
[20, 30]	02130	Male	Flu
[20, 30]	02140	Female	Ovarian Cancer
[20, 30]	02140	Male	Lung cancer
[30, 40]	02130	Female	Thyroid Disease
[30, 40]	02130	Male	Thyroid Disease



■ Bob

First Name	Surname	Age	Post Code	Gender
Alice	McCarthy	19	02138	Female
Steve	Joe	20	02139	Male
Eve	Boxwell	25	02141	Female
John	Mark	22	02142	Male
Rochelle	Alers	35	02138	Female
Ben	Colbert	30	02139	Male

Voter Registration Data



Unsorted Matching Attack...

First Name	Surname	Age	Post Code	Gender
Alice	McCarthy	19	02138	Female
Steve	Joe	20	02139	Male
Eve	Boxwell	25	02141	Female
John	Mark	22	02142	Male
Rochelle	Alers	35	02138	Female
Ben	Colbert	30	02139	Male

Records appear in the same order in the released table as in the original table

Age	Post Code	Gender	Symptoms
[20, 30]	02130	Female	Ovarian Cancer
[20, 30]	02130	Male	Flu
[20, 30]	02140	Female	Ovarian Cancer
[20, 30]	02140	Male	Lung cancer
[30, 40]	02130	Female	Thyroid Disease
[30, 40]	02130	Male	Thyroid Disease

Randomise before sharing? (Solution)

Complementary Release Attack...

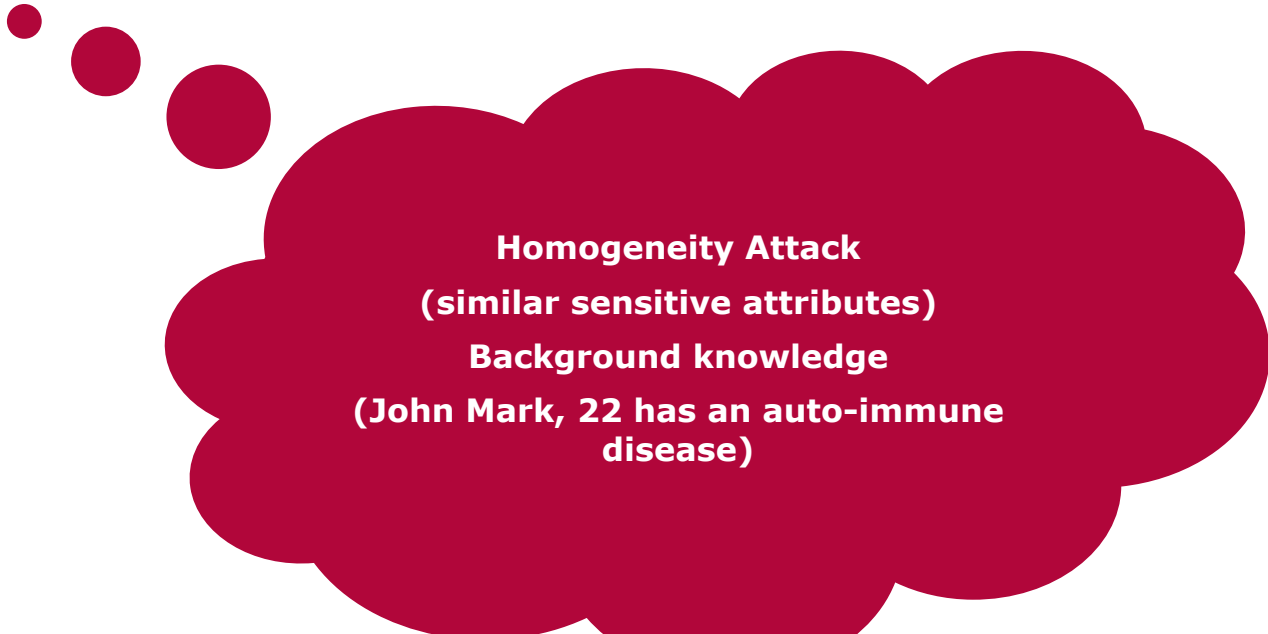
Age	Post Code	Gender	Symptoms
[20, 30]	02130	Female	Ovarian Cancer
[20, 30]	02130	Male	Flu
[20, 30]	02140	Female	Ovarian Cancer
[20, 30]	02140	Male	Lung cancer
[30, 40]	02130	Female	Thyroid Disease
[30, 40]	02130	Male	Thyroid Disease

Different releases of the same private table can be linked together to compromise k-anonymity

Race	Nationality	Age	Post Code	Gender
White	Belgian	[20, 30]	02138	Female
White	Australian	[20, 30]	0213*	Male
Mixed	Iceland	[20, 30]	02141	Female
Black	Canadian	[20, 30]	0214*	Male
Black	Jamaican	[30, 40]	02138	Female
Indian	Indian	[30, 40]	0213*	Male

Homogeneity and Background Knowledge...

Age	Post Code	Gender	Symptoms
[20, 30]	02130	Female	Heart Disease
[20, 30]	02130	Male	Heart Disease
[20, 30]	02140	Female	Heart Disease
[20, 30]	02140	Male	IBS
[30, 40]	02130	Female	Thyroid Disease
[30, 40]	02130	Male	Thyroid Disease



Similarity and Skewness Attack...

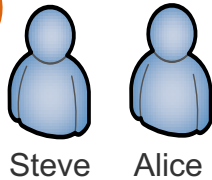
(Li et al. , 2007; Machanavajjhala et al., 2007)

- Utility versus Privacy

Age	Post Code	Gender	Symptoms
[20, 30]	02130	Female	Ovarian Cancer
[20, 30]	02130	Male	Flu
[20, 30]	02140	Female	Ovarian Cancer
[20, 30]	02140	Male	Lung Cancer
[30, 40]	02130	Female	Thyroid Disease
[30, 40]	02130	Male	Thyroid Disease

L-diversity does not consider overall distribution of sensitive values!

Similarity attack – exploits semantic closeness of attributes



Sensitive Attribute

- Machanavajjhala et al., “l-diversity: Privacy Beyond k-anonymity”, ACM Trans. On Knowl. Discov. Data 1(1), 2007
- Li et al., “t-Closeness: Privacy beyond k-anonymity and l-diversity”, IEEE 23rd International Conference on Data Engineering, pp. 106 – 115, 2007

Perfect k-anonymity is hard to get...

- Syntactic
 - Focuses on data transformation, not on what can be learned from the the anonymised dataset
 - “k-anonymous” dataset can leak sensitive information
 - “Quasi-identifier” fallacy
 - Assumes a priori that the attacker has no knowledge about the target
 - Relies on locality
 - Difficult to apply to real-world datasets (usually combined with other techniques)

**Do you think Network
Anonymisation is the answer?**

Why? And Why not?

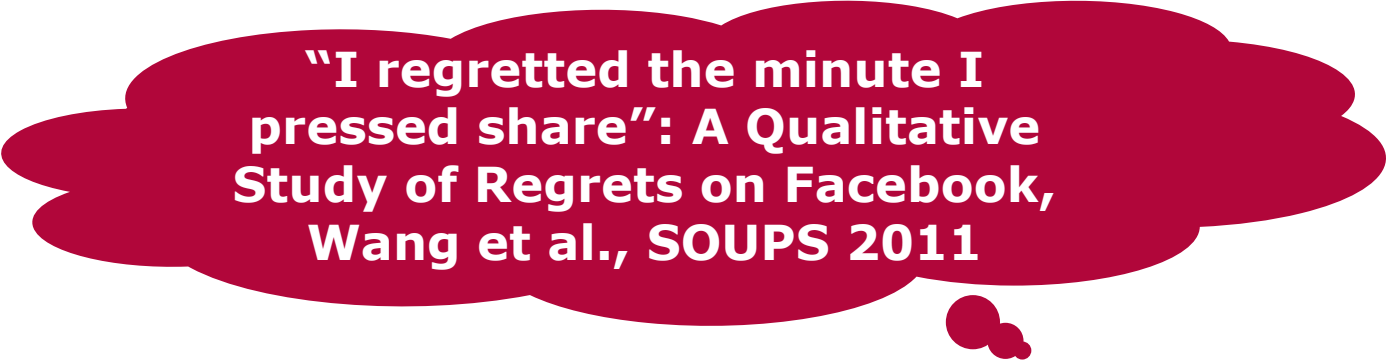
Next Week...

- Qualitative Research Methods
 - Surveys
 - Interviews
 - Diary Studies
 - Focus Groups
 - Crowdsourced Studies

Think About... (For next week...)

- See Handout (Notes2.pdf)
- Your proposal (for your project)
 - What would you like to do? Why?
 - Think of several options... make notes
 - What do you need? Initial usability data? Prototype?
 - Interesting adversarial scenarios?
 - Solutions to protect data?

Example of an Interview Study...



“I regretted the minute I pressed share”: A Qualitative Study of Regrets on Facebook, Wang et al., SOUPS 2011

- Methods Used: Reader Comments NYT, Survey, **Interviews**, and Diaries
- Pre-questionnaire for interviews:
 - Pittsburgh Craigslist
 - “Selected” 19 users from 301
 - Compensated \$20
 - Audio recorded and transcribed interviews along with screen shots

Example of an Interview Study...

Sample Interview Format...

Introduction: Say who you are, welcome the participant and tell them what the study is about...

Informed Consent: Describe the consent form and get consent from the user (signature)...

Interview Questions: (Motivations, Use of Facebook, Privacy Expectations, use of settings, regret experiences...)

Conclude: Thank the participant for their time and ask if he/she has any questions.